Lecture 8: spectra of deep net (Pennington et al. $18'^1$)

$$h^{(\ell)} = \phi(\hat{h}^{(\ell)})$$
 with $\hat{h}^{(\ell)} = W^{(\ell)} h^{(\ell-1)} + b^{(\ell)}$

has input to output Jacobian given by

$$J = \frac{\partial h^{(L)}}{\partial x^{(0)}} = \prod_{\ell=1}^{L} D^{(\ell)} W^{(\ell)}$$



Figure 4: Two limiting universality classes of Jacobian spectra. Hard Tanh and Shifted ReLU fall into one class, characterized by Bernoulli-distributed $\phi'(h)^2$, while Erf and Smoothed ReLU fall into a second class, characterized by a smooth distribution for $\phi'(h)^2$. The black curves are theoretical predictions for the limiting distributions with variance $\sigma_0^2 = 1/4$. The colored lines are emprical spectra of finite-depth width-1000 orthogonal neural networks. The empirical spectra converge to the limiting distributions in all cases. The rate of convergence is converge to the same limiting distribution along distinct fragietories. In all cases, the solid colored lines go from shallow L = 2 networks (red) to deep networks (purple). In all cases but Erf the deepest networks have L = 128. For Erf, the dashed lines show solutions to [15] for very large depth up to L = 8192.

¹https://arxiv.org/pdf/1802.09979.pdf

Theories of DL Lecture 9 Training a deep net: optimisation methods

- Backpropogation and the Pennington spectra
- Glorot (Xavier) initialisation to maintain variance through depth
- Batch normalization as a learned way to control saturation
- Methods for learning the net parameters:
 - Stochastic gradient descent (SGD)
 - Polyak and Nesterov momentum
 - AdaGrad and variants toward Adam

The output of the net is $H(x_{\mu}; \theta) = \hat{h}_{\mu}$ and we measure the value of the net through the average sum of squares:

$$\mathcal{L}(\theta; X, Y) = (2m)^{-1} \sum_{\mu=1}^{m} \sum_{i=1}^{n} (\hat{h}_{i,\mu} - y_{i,\mu})^2$$

Central to the success of deep nets is the ability to learn the parameters θ of the network to achieve good training error while also avoiding overfitting so as to generalize well.

Backpropogation allows an efficient calculation of $\operatorname{grad}_{\theta} L(\theta; X, Y)$.

$$\mathcal{L}(\theta; X, Y) = (2m)^{-1} \sum_{\mu=1}^{m} \sum_{i=1}^{n} (\hat{h}_{i,\mu} - y_{i,\mu})^2$$

Letting $\delta_{\ell} := \frac{\partial \mathcal{L}}{\partial \hat{h}^{(\ell)}}$ and as before $D^{(\ell)}$ the diagonal matrix with $D_{ii}^{(\ell)} = \phi'(\hat{h}_i^{(\ell)})$ we have

$$\delta_\ell = D^\ell (W^{(\ell)})^T \delta_{\ell+1}$$
 and $\delta_L = D^{(L)} \operatorname{grad}_{h^{(L)}} \mathcal{L}.$

which gives the formula for computing the δ_ℓ for each layer as

$$\delta_{\ell} = \left(\Pi_{k=\ell}^{L-1} D^{(k)} (W^{(k)})^T \right) D^{(L)} \operatorname{grad}_{h^{(L)}} \mathcal{L}.$$

and the resulting gradient $\text{grad}_{\theta}\mathcal{L}$ with entries as

$$\frac{\partial \mathcal{L}}{\partial W^{(\ell)}} = \delta_{\ell+1} \cdot h_{\ell}^{\mathsf{T}} \quad \text{ and } \frac{\partial \mathcal{L}}{\partial b^{(\ell)}} = \delta_{\ell+1}$$

Distribution of Jacobian spectra (Pennington et al. 18'²)

Recall the Jacobian of the input-output map

$$J = \frac{\partial h^{(L)}}{\partial x^{(0)}} = \prod_{\ell=1}^{L} D^{(\ell)} W^{(\ell)}$$

and contrast with gradient δ_ℓ

$$\delta_{\ell} = \left(\Pi_{k=\ell}^{L-1} D^{(k)} (W^{(k)})^T \right) D^{(L)} \operatorname{grad}_{h^{(L)}} \mathcal{L}.$$

we see the matrices with the same spectra and limiting distributions.



Figure 4: Two limiting universality classes of Jacobian spectra. Hard Tanh and Shifted ReLU fall into one class, characterized by Bernoulli-distributed $\phi'(h)^2$, while Erf and Smoothed ReLU fall into a second class, characterized by a smooth distribution for $\phi'(h)^2$. The black curves are theoretical predictions for the limiting distributions with variance $\sigma_0^2 = 1/4$. The colored lines are emprical spectra of finite-depth width-1000 orthogonal neural networks. The empirical spectra converge to the limiting distributions in all cases. The rate of convergence is converge to the same limiting distribution along distinct fragietories. In all cases, the solid colored lines go from shallow L = 2 networks (red) to deep networks (purple). In all cases but Erf the deepest networks have L = 128. Ever Erf, the dashed lines show solutions to CD. for vary large depth up to L = 8192.

²https://arxiv.org/pdf/1802.09979.pdf

Theories of DL Lecture 9 Training a deep net: optimisation methods

Backpropogation: weight initialization (Glorot et al.' 10^3)

An earlier, more empirical, perspective on the same issue of vanishing/exploding gradients was considered by Xavier Glorot and Yoshua Bengio (2010). They considered the same model assumptions as Pennington, \hat{h}^{ℓ} being approximately $\mathcal{N}(0, \sigma_{\ell}^2)$ and: "Our objective heres is to understand why standard gradient descent from random initialization is doing so poorly with deep neural networks.... Finally, we study how activations and gradients vary across layers and during training, with the idea that training may be more difficult when the singular values of the Jacobian associated with each layer are far from 1."



Figure 2: Mean and standard deviation (vertical bars) of the activation values (output of the sigmoid) during supervised learning, for the different hidden layers of a deep architecture. The top hidden layer quickly saturates at 0 (slowing down all learning), but then slowly desaturates around epoch 100.

³http://proceedings.mlr.press/v9/glorot10a.html

Backpropogation: weight initialization (Glorot et al.' 10^4)

Glorot initialization follows from seeking the variance of both the backpropogation gradient and forward activations to maintain the same variance per layer: for $W^{(\ell)} \in \mathbb{R}^{n \times n}$ need $\sigma_w^2 = 1/3n$.





Figure 6: Activation values normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized initialization (bottom). Top: 0-peak increases for higher layers.

Figure 7: Back-propagated gradients normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized (bottom) initialization. Top: 0-peak decreases for higher layers.



Figure 11: Test error during online training on the Shapeset-3 × 2 dataset. For various activation functions and heading fractions and the second statement of the second second

⁴http://proceedings.mlr.press/v9/glorot10a.html

Batch normalization (loffe et al. $15'^5$)

Saturation in a deep net can alternatively be controlled in a more learned manner by trying to shape the input to a nonlinear activation to be with a learned mean and variance:

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$; Parameters to be learned: γ, β **Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$ $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$ // mini-batch mean $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$ // mini-batch variance $\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ // normalize $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$ // scale and shift

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Typically included after a fully connected layer, before the nonlinear activation. Performed independently per nonlinear activation, with the γ and β learned as part of the net parameters θ .

⁵https://arxiv.org/pdf/1502.03167.pdf

Batch normalization experiment (loffe et al. 15'⁶)



Figure 2: Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.

Model	Steps to 72.2%	Max accuracy
Inception	$31.0 \cdot 10^{6}$	72.2%
BN-Baseline	$13.3\cdot 10^6$	72.7%
BN-x5	$2.1\cdot 10^6$	73.0%
BN-x30	$2.7\cdot 10^6$	74.8%
BN-x5-Sigmoid		69.8%

Figure 3: For Inception and the batch-normalized variants, the number of training steps required to reach the maximum accuracy of Inception (72.2%), and the maximum accuracy achieved by the network.

⁶https://arxiv.org/pdf/1502.03167.pdf

Stochastic gradient descent (SGD)

Given a loss function $\mathcal{L}(\theta; X, Y)$, gradient descent is given by $\theta^{(k+1)} = \theta^{(k)} - \eta \cdot \operatorname{grad}_{\theta} \mathcal{L}(\theta, X, Y)$

where η is referred to as the stepsize, or in deep learning the "learning rate."

Recall, we typically have a loss function which is the sum of *n* individual loss functions, independent for each data point: $\mathcal{L}(\theta; X, Y) = n^{-1} \sum_{\mu=1}^{n} l(\theta; x_{\mu}, y_{\mu})$ For $n \gg 1$ gradient descent is computationally too costly and instead one can break appart the *n* loss functions into "mini-batches" and repeatedly solve

$$\theta^{(k+1)} = \theta^{(k)} - \eta |\Lambda_k|^{-1} \operatorname{grad}_{\theta} \sum_{\mu \in \Lambda_k} I(\theta; x_{\mu}, y_{\mu}).$$

This is referred to as stochastic gradient descent as typically Λ_k is chosen in some randomized method, usually as a partition of [n] and a sequence of Λ_k which cover [n] is referred to as an "epoch."

Stochastic gradient descent: challenges and benefits

$$\theta^{(k+1)} = \theta^{(k)} - \eta |\Lambda_k|^{-1} \operatorname{grad}_{\theta} \sum_{\mu \in \Lambda_k} I(\theta; x_{\mu}, y_{\mu}).$$

- SGD is preferable for large n as it reduces the per iteration computational cost dependence on n to instead depend on |Λ_k| which can be set by the user as opposed to n which is given by the data set.
- SGD, and gradient descent, require selection of a learning rate (stepsize) which in deep learning is typically selected using some costly trial and error heuristics.
- The learning rate is typically chosen adaptively in a way that satisfies $\sum_{k=1}^{\infty} \eta_k = \infty$ and $\sum_{k=1}^{\infty} \eta_k^2 < \infty$; in particular as $\eta_k \sim k^{-1}$.
- The optimal selection of learning weight, and selection of Λ, depends on the unknown local Lipschitz constant ||grad/(θ₁; x_μ, y_μ) − grad/(θ₂; x_μ, y_μ)|| ≤ L_μ||θ₁ − θ₂||.

There are many improvements of stochastic gradient descent typically used in practise for deep learning; particularly popular is Polyak momentum:

$$\theta^{(k+1)} = \theta^{(k)} + \beta(\theta^{(k)} - \theta^{(k-1)}) - \alpha \cdot \operatorname{grad}_{\theta} \mathcal{L}\left(\theta^{(k)}\right)$$

or Nesterov's accelerated gradient:

$$\hat{\theta}^{k} = \theta^{(k)} + \beta(\theta^{(k)} - \theta^{(k-1)}) \theta^{(k+1)} = \hat{\theta}^{(k)} - \alpha \cdot \operatorname{grad}_{\theta} \mathcal{L}\left(\hat{\theta}^{(k)}\right)$$

These acceleration methods give substantial improvements in the linear convergence rate for convex problems; linear convergence rates are: Normal GD $\frac{\kappa-1}{\kappa+1}$, Polyak $\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$ and NAG $\sqrt{\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}}}$.